

## Laboratorium 3. Stemming, lematyzacja, worek słów, taka klasyfikacja

Poprzednie zajęcia pozwoliły poznać narzędzia do gromadzenia danych. Mając dane oznaczone etykietami, można wykonać prosty eksperyment klasyfikacji. Jednakże, surowy tekst napisany w języku naturalnym nie nadaje się do trenowania modeli. Aby było to możliwe konieczna jest wektoryzacja danych. Rozwiązanie tego laboratorium pozwoli poznać przykładowy sposób transformacji surowych danych tekstowych z etykietami do prostego modelu zdolnego do klasyfikacji.

### Zadanie 3.1. Przygotowanie danych

Wczytaj dane z katalogów `positive` oraz `negative` oraz dla każdego dokumentu:

1. Zamień w tekście wszystkie duże litery na małe
2. Dokonaj tokenizacji z użyciem `TreebankWordTokenizer`
3. Usuń z tokenów stopwords<sup>1</sup> (możesz też usunąć inne tokeny, które nie mają znaczenia np. "<br />")
4. Wykonaj stemming na tokenach z poprzedniego punktu z użyciem `PorterStemmer`
5. Wykonaj lematyzację na tokenach z punktu 3. z użyciem `WordNetLemmatizer`
6. Zapisz do tablic cztery zestawy tokenów:
  - Tokeny
  - Bez stop words
  - Po stemmingu
  - Po lematyzacji

#### Dokumentacja

- `TreebankWordTokenizer`<sup>2</sup>
- `PorterStemmer`<sup>3</sup>
- `WordNetLemmatizer`<sup>4</sup>

### Zadanie 3.2. Worek słów

Wykonaj wektoryzację dla każdego z czterech zestawów tokenów z zadania 3.1, zgodnie z podejściem "Bag of Words"

- Początkowo zapisz dane do zwykłej listy
- Wykorzystaj funkcję `Counter` w celu wyznaczenia występowania słów rozszerzając listę
- Zamień listę na macierz `DataFrame` wykorzystując funkcję `from_records`
- Zamień wartości "Nan" na zera oraz przekonwertuj na macierz typu `ndarray`<sup>5</sup>

#### Dokumentacja

- `Counter`<sup>6</sup>
- `from_records`<sup>7</sup>
- `fillna`<sup>8</sup>

### Zadanie 3.3. Klasyfikacja

1. Utwórz tablicę etykiet (300 zer oraz 300 jedynek) typu `ndarray`
2. Podziel dane na treningowe (70%) oraz testowe (30%) z użyciem `train_test_split`
3. Wytrenuj model z użyciem klasyfikatora `MultinomialNB` na danych treningowych (funkcja `fit()`)
4. Dokonaj predykcji modelu na danych testowych i zapisz wynik (funkcja `predict()`)
5. Wyświetl dokładności wyrażoną za pomocą metryki `accuracy_score`

#### Dokumentacja

- `train_test_split`<sup>9</sup>
- `MultinomialNB`<sup>10</sup>
- `accuracy_score`<sup>11</sup>

<sup>1</sup><https://www.nltk.org/book/ch02.html>

<sup>2</sup><https://www.nltk.org/api/nltk.tokenize.treebank.html>

<sup>3</sup><https://www.nltk.org/api/nltk.stem.porter.html#module-nltk.stem.porter>

<sup>4</sup><https://www.nltk.org/api/nltk.stem.wordnet.html>

<sup>5</sup><https://numpy.org/doc/stable/reference/generated/numpy.array.html>

<sup>6</sup><https://docs.python.org/3/library/collections.html#collections.Counter>

<sup>7</sup>[https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.from\\_records.html?highlight=from\\_records](https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.from_records.html?highlight=from_records)

<sup>8</sup><https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.fillna.html>

<sup>9</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.model\\_selection.train\\_test\\_split.html](https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.train_test_split.html)

<sup>10</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.naive\\_bayes.MultinomialNB.html](https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.MultinomialNB.html)

<sup>11</sup>[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.accuracy_score.html)